48

solid substrate are provided, e.g., in Merrifield (1986)
Science 232:341-347, which is hereby incorporated herein by
reference.

C.    Labeling Target Nucleotides
       The labeling procedures used in the sequencing
embodiments will also be applicable in the fingerprinting
embodiments.  However, since the fingerprinting embodiments
often will involve relatively large target molecules and
relatively short oligonucleotide probes, the amount of signal
necessary to incorporate into the target sequence may be less
critical than in the sequencing applications.  For example, a
relatively long target with a relatively small number of labels
per molecule may be easily amplified or detected because of the
relatively large target molecule size.
       In various embodiments, it may be desired to cleave
the target into smaller segments as in the sequencing
embodiments.  The labeling procedures and cleavage techniques
described in the sequencing embodiments would usually also be
applicable here.

D.    Hybridization Conditions
       The hybridization conditions used in fingerprinting
embodiments will typically be less critical than for the
sequencing embodiments.  The reason is that the amount of
mismatching which may be useful in providing the fingerprinting
information would typically be far greater than that necessary
in sequencing uses.  For example, Southern hybridizations do
not typically distinguish between slightly mismatched
sequences.  Under these circumstances, important and valuable
information may be arrived at with less stringent hybridization
conditions while providing valuable fingerprinting information.
 However, since the entire substrate is typically exposed to
the target molecule at one time, the binding affinity of the
probes should usually be of approximately comparable levels.
For this reason, if oligonucleotide probes are being used,
their lengths should be approximately comparable and will be
selected to hybridize under conditions which are common for

49

49

most of the probes on the substrate. Much as in a Southern hybridization, the target and oligonucleotide probes are of lengths typically greater than about 25 nucleotides. Under appropriate hybridization conditions, e.g., typically higher salt and lower temperature, the probes will hybridize irrespective of imperfect complementarity. In fact, with probes of greater than, e.g., about fifty nucleotides, the difference in stability of different sized probes will be relatively minor.

Typically the fingerprinting is merely for probing similarity or homology. Thus, the stringency of hybridization can usually be decreased to fairly low levels. See, e.g., Wetmur and Davidson (1968) "Kinetics of Renaturation of DNA," J. Mol. Biol., 31:349-370; and Kanehisa, M. (1984) Nuc. Acids Res., 12:203-213.

E.    Detection: VLSIPS™ Technology Scanning

Detection methods will be selected which are appropriate for the selected label. The scanning device need not necessarily be digitized or placed into a specific digital database, though such would most likely be done. For example, the analysis in fingerprinting could be photographic. Where a standardized fingerprint substrate matrix is used, the pattern of hybridizations may be spatially unique and may be compared photographically. In this manner, each sample may have a characteristic pattern of interactions and the likelihood of identical patterns will preferably be such low frequency that the fingerprint pattern indeed becomes a characteristic pattern virtually as unique as an individual's fingertip fingerprint. With a standardized substrate, every individual could be, in theory, uniquely identifiable on the basis of the pattern of hybridizing to the substrate.

Of course, the VLSIPS™ Technology scanning apparatus may also be useful to generate a digitized version of the fingerprint pattern. In this way, the identification pattern can be provided in a linear string of digits. This sequence could also be used for a standardized identification system providing significant useful medical transferability of

50

50

specific data.  In one embodiment, the probes used are selected
to be of sufficiently high resolution to measure the antigens
of the major histo compatibility complex.  It might even be
possible to provide transplantation matching data in a linear
stream of data.  The fingerprinting data may provide a
condensed version, or summary, of the linear genetic data, or
any other information data base.

F.   Analysis

The analysis of the fingerprint will often be much
simpler than a total sequence determination.  However, there
may be particular types of analysis which will be substantially
simplified by a selected group of probes.  For example, probes
which exhibit particular populational heterogeneity may be
selected.  In this way, analysis may be simplified and
practical utility enhanced merely by careful selection of the
specific probes and a careful matrix layout of those probes.

G.   Substrate Reuse

As with the sequencing application, the
fingerprinting usages may also take advantage of the
reusability of the substrate.  In this way, the interactions
can be disrupted, the substrate treated, and the renewed
substrate is equivalent to an unused substrate.

H.   Non-polynucleotide Aspects

Besides polynucleotide applications, the
fingerprinting analysis may be applied to other polymers,
especially polypeptides, carbohydrates, and other polymers,
both organic and inorganic.  Besides using the fingerprinting
method for analyzing a particular polymer, the fingerprinting
method may be used to characterize various samples.  For
example, a cell or population of cells may be tested for their
expression of specific antigens or their mRNA sequence intent.
For example, a T-cell may be classified by virtue of its
combination of expressed surface antigens.  With specific
reagents which interact with these antigens, a cell or a
population of cells or a lysed cell may be exposed to a VLSIPS

51

51

substrate. The biological sample may be classified or characterized by analyzing the pattern of specific interaction. This may be applicable to a cell or tissue type, to the messenger RNA population expressed by a cell to the genetic content of a cell, or to virtually any sample which can be classified and/or identified by its combination of specific molecular properties.

The ability to generate a high density means for screening the presence or absence of specific interactions allows for the possibility of screening for, if not saturating, all of a very large number of possible interactions. This is very powerful in providing the means for testing the combinations of molecular properties which can define a class of samples. For example, a species of organism may be characterized by its DNA sequences, e.g., a genetic fingerprint. By using a fingerprinting method, it may be determined that all members of that species are sufficiently similar in specific sequences that they can be easily identified as being within a particular group. Thus, newly defined classes may be resolved by their similarity in fingerprint patterns. Alternatively, a non-member of that group will fail to share those many identifying characteristics. However, since the technology allows testing of a very large number of specific interactions, it also provides the ability to more finely distinguish between closely related different cells or samples. This will have important applications in diagnosing viral, bacterial, and other pathological on nonpathological infections.

In particular, cell classification may be defined by any of a number of different properties. For example, a cell class may be defined by its DNA sequences contained therein. This allows species identification for parasitic or other infections. For example, the human cell is presumably genetically distinguishable from a monkey cell, but different human cells will share many genetic markers. At higher resolution, each individual human genome will exhibit unique sequences that can define it as a single individual.

Likewise, a developmental stage of a cell type may be

52

definable by its pattern of expression of messenger RNA. For example, in particular stages of cells, high levels of ribosomal RNA are found whereas relatively low levels of other types of messenger RNAs may be found. The high resolution distinguishability provided by this fingerprinting method allows the distinction between cells which have relatively minor differences in its expressed mRNA population. Where a pattern is shown to be characteristic of a stage, a stage may be defined by that particular pattern of messenger RNA expression.

In a similar manner, the antigenic determinants found on a protein may very well define the cell class. For example, immunological T-cells are distinguishable from B-cells because, in part, the cell surface antigens on the cell types are distinguishable. Different T-cell subclasses can be also distinguished from one another by whether they contain particular T-cell antigens. The present invention provides the possibility for high resolution testing of many different interactions simultaneously, and the definition of new cell types will be possible.

The high resolution VLSIPS™ substrate may also be used as a very powerful diagnostic tool to test the combination of presence, of a plurality of different assays from a biological sample. For example, a cancerous condition may be indicated by a combination of various different properties found in the blood. For example, a cancerous condition may be indicated by a combination of expression of various soluble antigens found in the blood along with a high number of various cellular antigens found on lymphocytes and/or particular cell degradation products. With a substrate as provided herein, a large number of different features can be simultaneously performed on a biological sample. In fact, the high resolution of the test will allow more complete characterization of parameters which define particular diseases. Thus, the power of diagnostic tests may be limited by the extent of statistical correlation with a particular condition rather than with the number of antigens or interactions which are tested. The present invention provides the means to generate this large

53

53

universe of possible reagents and the ability to actually
accumulate that correlative data.

In another embodiment, a substrate as provided herein
may be used for genetic screening. This would allow for
simultaneous screening of thousands of genetic markers. As the
density of the matrix is increased, many more molecules can be
simultaneously tested. Genetic screening then becomes a
simpler method as the present invention provides the ability to
screen for thousands, tens of thousands, and hundreds of
thousands, even millions of different possible genetic
features. However, the number of high correlation genetic
markers for conditions numbers only in the hundreds. Again,
the possibility for screening a large number of sequences
provides the opportunity for generating the data which can
provide correlation between sequences and specific conditions
or susceptibility. The present invention provides the means to
generate extremely valuable correlations useful for the genetic
detection of the causative mutation leading to medical
conditions. In still another embodiment, the present invention
would be applicable to distinguishing two individuals having
identical genetic compositions. The antibody population within
an individual is dependent both on genetic and historical
factors. Each individual experiences a unique exposure to
various infectious agents, and the combined antibody expression
is partly determined thereby. Thus, individuals may also be
fingerprinted by their immunological content, either of
actively expressed antibodies, or their immunological memory.
Similar sorts of immunological and environmental histories may
be useful for fingerprinting, perhaps in combination with other
screening properties. In particular, the present invention may
be useful for screening allergic reactions or susceptibilities,
and a simple IgE specificity test may be useful in determining
a spectrum of allergies.

With the definition of new classes of cells, a cell
sorter will be used to purify them. Moreover, new markers for
defining that class of cells will be identified. For example,
where the class is defined by its RNA content, cells may be
screened by antisense probes which detect the presence or

54

54

absence of specific sequences therein.  Alternatively, cell
lysates may provide information useful in correlating
intracellular properties with extracellular markers which
indicate functional differences.  Using standard cell sorter
technology with a fluorescence or labeled antisense probe which
recognizes the internal presence of the specific sequences of
interest, the cell sorter will be able to isolate a relatively
homogeneous population of cells possessing the particular
marker.  Using successive probes the sorting process should be
able to select for cells having a combination of a large number
of different markers.

In a non-polynucleotide embodiment, cells may be
defined by the presence of other markers.  The markers may be
carbohydrates, proteins, or other molecules.  Thus, a substrate
having particular specific reagents, e.g., antibodies, attached
to it should be able to identify cells having particular
patterns of marker expression.  Of course, combinations of
these made be utilized and a cell class may be defined by a
combination of its expressed mRNA, its carbohydrate expression,
its antigens, and other properties.  This fingerprinting should
be useful in determining the physiological state of a cell or
population of cells.

Having defined a cell type whose function or
properties are defined by the reagents attachable to a VLSIPS
substrate, such as cellular antigens, these structural
manifestations of function may be used to sort cells to
generate a relatively homogeneous population of that class of
cells.  Standard cell sorter technology may be applied to
purify such a population, see, e.g., Dangl, J. and Herzenberg
(1982) "Selection of hybridomas and hybridoma variants using
the fluorescence activated cell sorter," J. Immunological
Methods 52:1-14; and Becton Dickinson, Fluorescence Activated
Cell Sorter Division, San Jose, California, and Coulter
Diagnostics, Hialeah, Florida.

With the fingerprinting method an identification
means arises from mosaicism problems in an organism.  A mosaic
organism is one whose genetic content in different cells is
significantly different.  Various clonal populations should

55

have similar genetic fingerprints, though different clonal
populations may have different genetic contents.  See, for
example, Suzuki et al. An Introduction to Genetic Analysis (4th
Ed.), Freeman and Co., New York, which is hereby incorporated
herein by reference.  However, this problem should be a
relatively rare problem and could be more carefully evaluated
with greater experience using the fingerprinting methods.

The invention will also find use in detecting
changes, both genetic and antigenic, e.g., in a rapidly
"evolving" protozoa infection, or similarly changing organism.

V.    MAPPING

A.    General

The use of the present invention for mapping
parallels its use for fingerprinting and sequencing.  Where a
polymer is a linear molecule, the mapping provides the ability
to locate particular segments along the length of the polymer.
Branched polymers can be treated as a series of individual
linear polymers.  The mapping provides the ability to locate,
in a relative sense, the order of various subsequences.  This
may be achieved using at least two different approaches.

The first approach is to take the large sequence and
fragment it at specific points.  The fragments are then ordered
and attached to a solid substrate.  For example, the clones
resulting from a chromosome walking process may be individually
attached to the substrate by methods, e.g., caged biotin
techniques, indicated earlier.  Segments of unknown map
position will be exposed to the substrate and will hybridize to
the segment which contains that particular sequence.  This
procedure allows the rapid determination of a number of
different labeled segments, each mapping requiring only a
single hybridization step once the substrate is generated.  The
substrate may be regenerated by removal of the interaction, and
the next mapping segment applied.

In an alternative method, a plurality of subsequences
can be attached to a substrate.  Various short probes may be
applied to determine which segments may contain particular
overlaps.  The theoretical basis and a description of this

56

mapping procedure is contained in, e.g., Evans et al. 1989
"Physical Mapping of Complex Genomes by Cosmid Multiplex
Analysis," Proc. Natl. Acad. Sci. USA 86:5030-5034, and other
references cited above in the Section labeled "Overall
Description." Using this approach, the details of the mapping
embodiment are very similar to those used in the fingerprinting
embodiment.

B.    Preparation of Substrate Matrix
     The substrate may be generated in either of the
methods generally applicable in the sequencing and
fingerprinting embodiments. The substrate may be made either
synthetically, or by attaching otherwise purified probes or
sequences to the matrix. The probes or sequences may be
derived either from synthetic or biological means. As
indicated above, the solid phase substrate synthetic methods
may be utilized to generate a matrix with positionally defined
sequences. In the mapping embodiment, the importance of
saturation of all possible subsequences of a preselected length
is far less important than in the sequencing embodiment, but
the length of the probes used may be desired to be much longer.
 The processes for making a substrate which has longer
oligonucleotide probes should not be significantly different
from those described for the sequencing embodiments, but the
optimization parameters may be modified to comply with the
mapping needs.

C.    Labeling
     The labeling methods will be similar to those
applicable in sequencing and fingerprinting embodiments.
Again, it may be desirable to fragment the target sequences.

D.    Hybridization/Specific Interaction
     The specificity of interaction between the targets
and probe would typically be closer to those used for
fingerprinting embodiments, where homology is more important
than absolute distinguishability of high fidelity complementary
hybridization. Usually, the hybridization conditions will be

57

such that merely homologous segments will interact and provide a positive signal. Much like the fingerprinting embodiment, it may be useful to measure the extent of homology by successive incubations at higher stringency conditions. Or, a plurality of different probes, each having various levels of homology may be used. In either way, the spectrum of homologies can be measured.

Where non-nucleic acid hybridization is involved, the specific interactions may also be compared in a fingerprint-like manner. The specific reagents may have less specificity, e.g., monoclonal antibodies which recognize a broader spectrum of sequences may be utilized relative to a sequencing embodiment. Again, the specificity of interaction may be measured under various conditions of increasing stringency to determine the spectrum of matching across the specific probes selected, or a number of different stringency reagents may be included to indicate the binding affinity.

E.    Detection

The detection methods used in the mapping procedure will be virtually identical to those used in the fingerprinting embodiment. The detection methods will be selected in combination with the labeling methods.

F.    Analysis

The analysis of the data in a mapping embodiment will typically be somewhat different from that in fingerprinting. The fingerprinting embodiment will test for the presence or absence of specific or homologous segments. However, in the mapping embodiment, the existence of an interaction is coupled with some indication of the location of the interaction. The interaction is mapped in some manner to the physical polymer sequence. Some means for determining the relative positions of different probes is performed. This may be achieved by synthesis of the substrate in pattern, or may result from analysis of sequences after they have been attached to the substrate.

For example, the probes may be randomly positioned at

58

58

various locations on the substrate. However, the relative
positions of the various reagents in the original polymer may
be determined by using short fragments, e.g., individually, as
target molecules which determine the proximity of different
probes. By an automated system of testing each different short
fragment of the original polymer, coupled with proper analysis,
it will be possible to determine which probes are adjacent one
another on the original target sequence and correlate that with
positions on the matrix. In this way, the matrix is useful for
determining the relative locations of various new segments in
the original target molecule. This sort of analysis is
described in Evans, and the related references described above.

G.    Substrate Reuse
The substrate should be reusable in the manner
described in the fingerprinting section. The substrate is
renewed by removal of the specific interactions and is washed
and prepared for successive cycles of exposure to new target
sequences.

H.    Non-polynucleotide Aspects
The mapping procedure may be used on other molecules
than polynucleotides. Although hybridization is one type of
specific interaction which is clearly useful for use in this
mapping embodiment, antibody reagents may also be very useful.
In the same way that polypeptide sequencing or other polymers
may be sequenced by the reagents and techniques described in
the sequencing section and fingerprinting section, the mapping
embodiment may also be used similarly.
In another form of mapping, as described above in the
fingerprinting section, the developmental map of a cell or
biological system may be measured using fingerprinting type
technology. Thus, the mapping may be along a temporal
dimension rather than along a polymer dimension. The mapping
or fingerprinting embodiments may also be used in determining
the genetic rearrangements which may be genetically important,
as in lymphocyte and B-cell development. In another example,
various rearrangements or chromosomal dislocations may be

59

tested by either the fingerprinting or mapping methods. These
techniques are similar in many respects and the fingerprinting
and mapping embodiments may overlap in many respects.

VI.    ADDITIONAL SCREENING AND APPLICATIONS
    A.    Specific Interactions
        As originally indicated in the parent filing of
VLSIPS™ Technology, the production of a high density plurality
of spatially segregated polymers provides the ability to
generate a very large universe or repertoire of individually
and distinct sequence possibilities. As indicated above,
particular oligonucleotides may be synthesized in automated
fashion at specific locations on a matrix. In fact, these
oligonucleotides may be used to direct other molecules to
specific locations by linking specific oligonucleotides to
other reagents which are in batch exposed to the matrix and
hybridized in a complementary fashion to only those locations
where the complementary oligonucleotide has been synthesized on
the matrix. This allows for spatially attaching a plurality of
different reagents onto the matrix instead of individually
attaching each separate reagent at each specific location.
Although the caged biotin method allows automated attachment,
the speed of the caged biotin attachment process is relatively
slow and requires a separate reaction for each reagent being
attached. By use of the oligonucleotide method, the
specificity of position can be done in an automated and
parallel fashion. As each reagent is produced, instead of
directly attaching each reagent at each desired position, the
reagent may be attached to a specific desired complementary
oligonucleotide which will ultimately be specifically directed
toward locations on the matrix having a complementary
oligonucleotide attached thereat.

        In addition, the technology allows screening for
specificity of interaction with particular reagents. For
example, the oligonucleotide sequence specificity of binding of
a potential reagent may be tested by presenting to the reagent
all of the possible subsequences available for binding.
Although secondary or higher order sequence specific features

60

60

might not be easily screenable using this technology, it does
provide a convenient, simple, quick, and thorough screen of
interactions between a reagent and its target recognition
sequences. See, e.g., Pfeifer et al. (1989) Science 246:810-
812.

For example, the interaction of a promoter protein
with its target binding sequence may be tested for many
different, or all, possible binding sequences. By testing the
strength of interactions under various different conditions,
the interaction of the promoter protein with each of the
different potential binding sites may be analyzed. The
spectrum of strength of interactions with each different
potential binding site may provide significant insight into the
types of features which are important in determining
specificity.

An additional example of a sequence specific
interaction between reagents is the testing of binding of a
double stranded nucleic acid structure with a single stranded
oligonucleotide. Often, a triple stranded structure is
produced which has significant aspects of sequence specificity.
Testing of such interactions with either sequences comprising
only natural nucleotides, or perhaps the testing of nucleotide
analogs may be very important in screening for particularly
useful diagnostic or therapeutic reagents. See, e.g., Häner
and Dervan (1990) Biochemistry 29:9761-6765, and references
therein.

### B.    Sequence Comparisons

Once a gene is sequenced, the present invention
provides a means to compare alleles or related sequences to
locate and identify differences from the control sequence.
This would be extremely useful in further analysis of genetic
variability at a specific gene locus.

### C.    Categorizations

As indicated above in the fingerprinting and mapping
embodiments, the present invention is also useful in defining
specific stages in the temporal sequence of cells, e.g.,

61

61

development, and the resulting tissues within an organism. For example, the developmental stage of a cell, or population of cells, can be dependent upon the expression of particular messenger RNAs or cellular antigens. The screening procedures provided allow for high resolution definition of new classes of cells. In addition, the temporal development of particular cells will be characterized by the presence or expression of various mRNAs. Means to simultaneously screen a plurality or very large number of different sequences are provided. The combination of different markers made available dramatically increases the ability to distinguish fairly closely related cell types. Other markers may be combined with markers and methods made available herein to define new classifications of biological samples, e.g., based upon new combinations of markers.

The presence or absence of particular marker sequences will be used to define temporal developmental stages. Once the stages are defined, fairly simple methods can be applied to actually purify those particular cells. For example, antisense probes or recognition reagents may be used with a cell sorter to select those cells containing or expressing the critical markers. Alternatively, the expression of those sequences may result in specific antigens which may also be used in defining cell classes and sorting those cells away from others. In this way, for example, it should be possible to select a class of omnipotent immune system cells which are able to completely regenerate a human immune system. Based upon the cellular classes defined by the parameters made available by this technology, purified classes of cells having identifiable differences, structural or functional, are made available.

In an alternative embodiment, a plurality of antigens or specific binding proteins attached to the substrate may be used to define particular cell types. For example, subclasses of T-cells are defined, in part, by the combination of expressed cell surface antigens. The present invention allows for the simultaneous screening of a large plurality of different antigens together. Thus, higher resolution

G⅄

62

classification of different T-cell subclasses becomes possible and, with the definitions and functional differences which correlate with those antigenic or other parameters, the ability to purify those cell types becomes available. This is applicable not only to T-cells, but also to lymphocyte cells, or even to freely circulating cells. Many of the cells for which this would be most useful will be immobile cells found in particular tissues or organs. Tumor cells will be diagnosed or detected using these fingerprinting techniques. Coupled with a temporal change in structure, developmental classes may also be selected and defined using these technologies. The present invention also provides the ability not only to define new classes of cells based upon functional or structural differences, but it also provides the ability to select or purify populations of cells which share these particular properties. Standard cell sorting procedures using antibody markers may be used to detect extracellular features. Intracellular features would also be detectable by introducing the label reagents into the cell. In particular, antisense DNA or RNA molecules may be introduced into a cell to detect RNA sequences therein. See, e.g., Weintraub (1990) Scientific American 262:40-46.

D.    Statistical Correlations

In an additional embodiment, the present invention also allows for the high resolution correlation of medical conditions with various different markers. For example, the presently available technology, when applied to amniocentesis or other genetic screening methods, typically screens for tens of different markers at most. The present invention allows simultaneous screening for tens, hundreds, thousands, tens of thousands, hundreds of thousands, and even millions of different genetic sequences. Thus, applying the fingerprinting methods of the present invention to a sufficiently large population allows detailed statistical analysis to be made, thereby correlating particular medical conditions with particular markers, typically antigenic or genetic. Tumor specific antigens will be identified using the present

63

63

invention.

Various medical conditions may be correlated against an enormous data base of the sequences within an individual. Genetic propensities and correlations then become available and high resolution genetic predictability and correlation become much more easily performed. With the enormous data base, the reliability of the predictions is also better tested. Particular markers which are partially diagnostic of particular medical conditions or medical susceptibilities will be identified and provide direction in further studies and more careful analysis of the markers involved. Of course, as indicated above in the sequencing embodiment, the present invention will find much use in intense sequencing projects. For example, sequencing of the entire human genome in the human genome project will be greatly simplified and enabled by the present invention.

VI.    FORMATION OF SUBSTRATE

The substrate is provided with a pattern of specific reagents which are positionally localized on the surface of the substrate. This matrix of positions is defined by the automated system which produces the substrate. The instrument will typically be one similar to that described in Pirrung et al. (1992) U.S. Pat. No. 5,143,854, and Serial No. 07/624,120, now abandoned. The instrumentation described therein is directly applicable to the applications used here. In particular, the apparatus comprises a substrate, typically a silicon containing substrate, on which positions on the surface may be defined by a coordinate system of positions. These positions can be individually addressed or detected by the VLSIPS™ Technology apparatus.

Typically, the VLSIPS™ Technology apparatus uses optical methods used in semiconductor fabrication applications. In this way, masks may be used to photo-activate positions for attachment or synthesis of specific sequences on the substrate. These manipulations may be automated by the types of apparatus described in Pirrung et al. (1992) U.S. Pat. No. 5,143,854 and Serial No. 07/624,120, now abandoned.

64

64

Selectively removable protecting groups allow creation of well defined areas of substrate surface having differing reactivities. Preferably, the protecting groups are selectively removed from the surface by applying a specific activator, such as electromagnetic radiation of a specific wavelength and intensity. More preferably, the specific activator exposes selected areas of surface to remove the protecting groups in the exposed areas.

Protecting groups of the present invention are used in conjunction with solid phase oligomer syntheses, such as peptide syntheses using natural or unnatural amino acids, nucleotide syntheses using deoxyribonucleic and ribonucleic acids, oligosaccharide syntheses, and the like. In addition to protecting the substrate surface from unwanted reaction, the protecting groups block a reactive end of the monomer to prevent self-polymerization. For instance, attachment of a protecting group to the amino terminus of an activated amino acid, such as the N-hydroxysuccinimide-activated ester of the amino acid prevents the amino terminus of one monomer from reacting with the activated ester portion of another during peptide synthesis.

Alternatively, the protecting group may be attached to the carboxyl group of an amino acid to prevent reaction at this site. Most protecting groups can be attached to either the amino or the carboxyl group of an amino acid, and the nature of the chemical synthesis will dictate which reactive group will require a protecting group. Analogously, attachment of a protecting group to the 5'-hydroxyl group of a nucleoside during synthesis using for example, phosphate-triester coupling chemistry, prevents the 5'-hydroxyl of one nucleoside from reacting with the 3'-activated phosphate-triester of another.

Regardless of the specific use, protecting groups are employed to protect a moiety on a molecule from reacting with another reagent. Protecting groups of the present invention have the following characteristics: they prevent selected reagents from modifying the group to which they are attached; they are stable (that is, they remain attached) to the synthesis reaction conditions; they are removable under

65

65

conditions that do not adversely affect the remaining
structure; and once removed, do not react appreciably with the
surface or surface-bound oligomer.  The selection of a suitable
protecting group will depend, of course, on the chemical nature
of the monomer unit and oligomer, as well as the specific
reagents they are to protect against.

In a preferred embodiment, the protecting groups will
be photoactivatable.  The properties and uses of photoreactive
protecting compounds have been reviewed.  See, McCray et al.,
Ann. Rev. of Biophys. and Biophys. Chem. (1989) 18:239-270,
which is incorporated herein by reference.  Preferably, the
photosensitive protecting groups will be removable by radiation
in the ultraviolet (UV) or visible portion of the
electromagnetic spectrum.  More preferably, the protecting
groups will be removable by radiation in the near UV or visible
portion of the spectrum.  In some embodiments, however,
activation may be performed by other methods such as localized
heating, electron beam lithography, laser pumping, oxidation or
reduction with microelectrodes, and the like.  Sulfonyl
compounds are suitable reactive groups for electron beam
lithography.  Oxidative or reductive removal is accomplished by
exposure of the protecting group to an electric current source,
preferably using microelectrodes directed to the predefined
regions of the surface which are desired for activation.  A
more detailed description of these protective groups is
provided in Serial No. 07/624,120, now abandoned, which is
hereby incorporated herein by reference.

The density of reagents attached to a silicon
substrate may be varied by standard procedures.  The surface
area for attachment of reagents may be increased by modifying
the silicon surface.  For example, a matte surface may be
machined or etched on the substrate to provide more sites for
attachment of the particular reagents.  Another way to increase
the density of reagent binding sites is to increase the
derivitization density of the silicon.  Standard procedures for
achieving this are described, below.

One method to control the derivatization density is
to highly derivatize the substrate with photochemical groups at

66

66

high·density.  The substrate is then photolyzed for various
predetermined times, which photoactivate the groups at a
measurable rate, and react them with a capping reagent.  By
this method, the density of linker groups may be modulated by
using a desired time and intensity of photoactivation.

In many applications, the number of different
sequences which may be provided may be limited by the density
and the size of the substrate on which the matrix pattern is
generated.  In situations where the density is insufficiently
high to allow the screening of the desired number of sequences,
multiple substrates may be used to increase the number of
sequences tested.  Thus, the number of sequences tested may be
increased by using a plurality of different substrates.
Because the VLSIPS apparatus is almost fully automated,
increasing the number of substrates does not lead to a
significant increase in the number of manipulations which must
be performed by humans.  This again leads to greater
reproducibility and speed in the handling of these multiple
substrates.

A.    Instrumentation
The concept of using VLSIPS™ Technology generally
allows a pattern or a matrix of reagents to be generated.  The
procedure for making the pattern is performed by any of a
number of different methods.  An apparatus and instrumentation
useful for generating a high density VLSIPS substrate is
described in detail in Pirrung et al. (1992) U.S. Pat. No.
5,143,854 and Serial No. 07/624,120, now abandoned.

.B.    Binary Masking
The details of the binary masking are described in an
accompanying application filed simultaneously with this, Serial
No. 07/624,120, now abandoned, whose specification is
incorporated herein by reference.

For example, the binary masking technique allows for
producing a plurality of sequences based on the selection of
either of two possibilities at any particular location.  By a
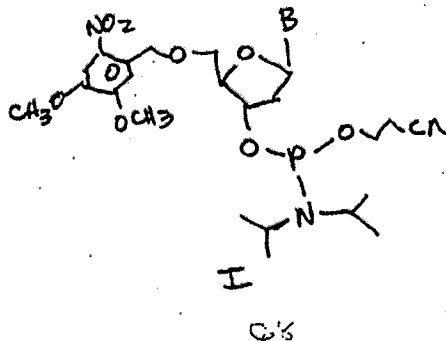series of binary masking steps, the binary decision may be the
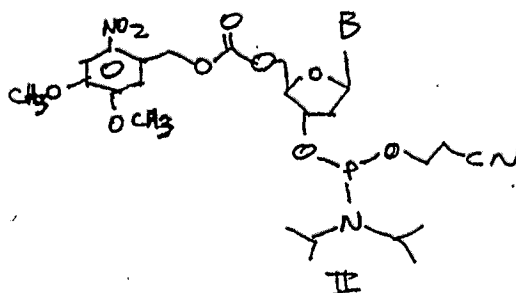
67

67

determination, on a particular synthetic cycle, whether or not
to add any particular one of the possible subunits.  By
treating various regions of the matrix pattern in parallel, the
binary masking strategy provides the ability to carry out
spatially addressable parallel synthesis.

C.    Synthetic Methods

The synthetic methods in making a substrate are
described in the parent application, Pirrung et al. (1992) U.S.
Pat. No. 5,143,854.  The construction of the matrix pattern on
the substrate will typically be generated by the use of photo-
sensitive reagents.  By use of photo-lithographic optical
methods, particular segments of the substrate can be irradiated
with light to activate or deactivate blocking agents, e.g., to
protect or deprotect particular chemical groups.  By an
appropriate sequence of photo-exposure steps at appropriate
times with appropriate masks and with appropriate reagents, the
substrates can have known polymers synthesized at positionally
defined regions on the substrate.  Methods for synthesizing
various substrates are described in Pirrung et al. (1992) U.S.
Pat. No. 5,143,854 and Serial No. 07/624,120, now abandoned.
By a sequential series of these photo-exposure and reaction
manipulations, a defined matrix pattern of known sequences may
be generated, and is typically referred to as a VLSIPS™
Technology substrate.  In the nucleic acid synthesis
embodiment, nucleosides used in the synthesis of DNA by
photolytic methods will typically be one of the two forms shown
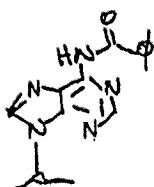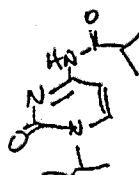below:
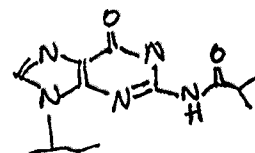
68



II

B = Adenine, Cytosine, Guanine, or Thymine

In I, the photolabile group at the 5' position is abbreviated NV (nitroveratryl) and in II, the group is abbreviated NVOC (nitroveratryl oxycarbonyl). Although not shown in Fig. C, the bases (adenine, cytosine, and guanine) contain exocyclic $NH_2$ groups which must be protected during DNA synthesis. Thymine contains no exocyclic $NH_2$ and therefore requires no protection. The standard protecting groups for these amines are shown below:



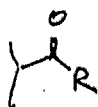Adenine (A)                 Cytosine (C)              Guanine (G)
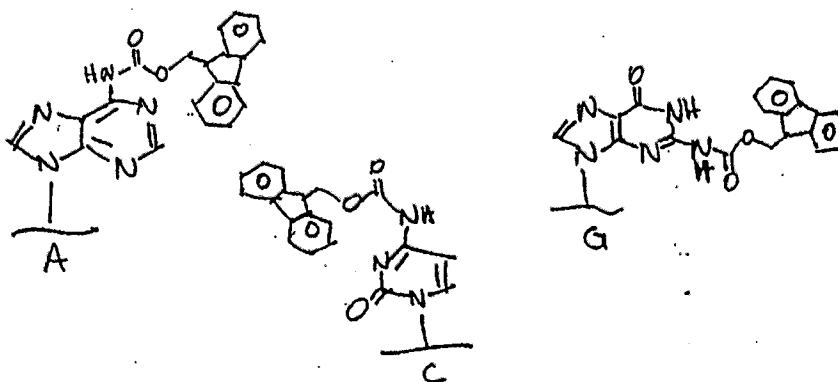
Other amides of the general formula

G.

69



where R may be alkyl or aryl have been used.

Another type of protecting group FMOC (9-fluorenyl methoxycarbonyl) is currently being used to protect the exocyclic amines of the three bases:



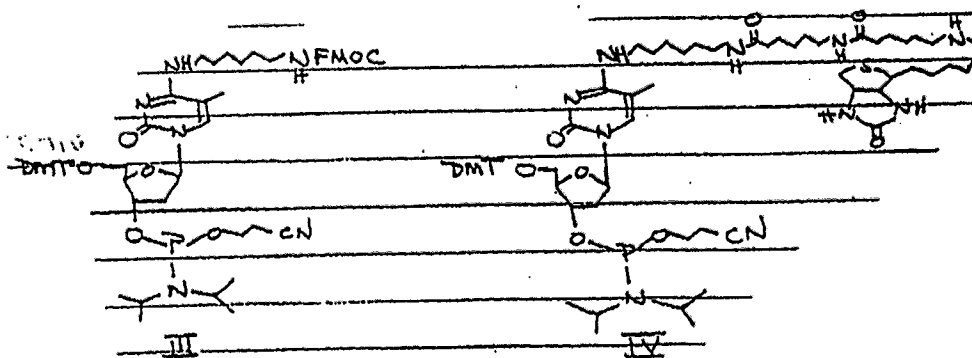Adenine (A)          Cytosine (C)          Guanine (G)

The advantage of the FMOC group is that it is removed under mild conditions (dilute organic bases) and can be used for all three bases. The amide protecting groups require more harsh conditions to be removed (NH,/MeOH with heat).
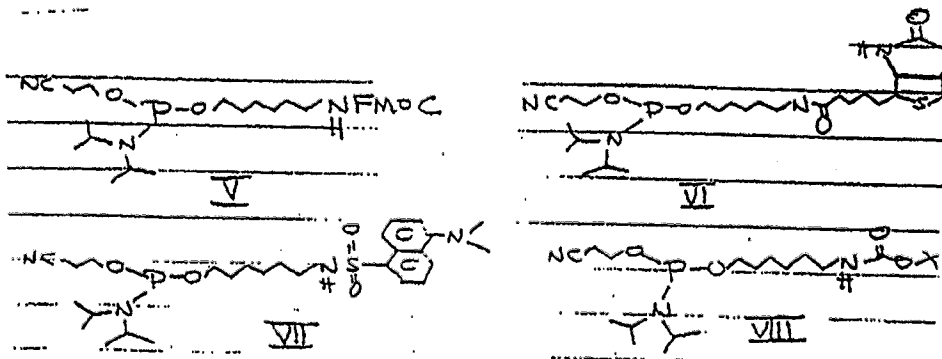
70

70

Nucleosides used as 5'-OH probes, useful in verifying
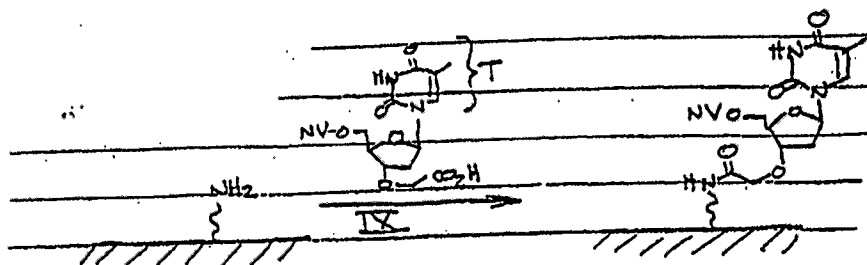correct VLSIPS synthetic function, include, for example, the
following:



These compounds are used to detect where on a
substrate photolysis has occurred by the attachment of either
III or V to the newly generated 5'-OH.  In the case of III,
after the phosphate attachment is made, the substrate is
treated with a dilute base to remove the FMOC group.  The
resulting amine can be reacted with FITC and the substrate
examined by fluorescence microscopy.  This indicates the proper
generation of a 5'-OH.  In the case of compound IV, after the
phosphate attachment is made, the substrate is treated with
FITC labeled streptavidin and the substrate again may be
examined by fluorescence microscopy.  Other probes, although
not nucleoside based, have included the following:
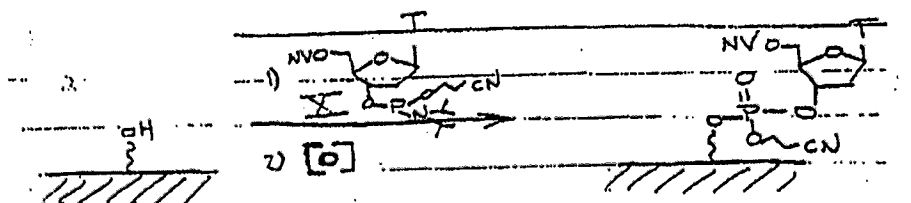


71

IAFP00000999

71

The method of attachment of the first nucleoside to
the surface of the substrate depends on the functionality of
the groups at the substrate surface.  If the surface is amine
functionalized, an amide bond is made (see example below).



If the surface is hydroxy functionalized, a phosphate
bond is made (see example below):



In both cases, the thymidine example is illustrated,
but any one of the four phosphoramidite activated nucleosides
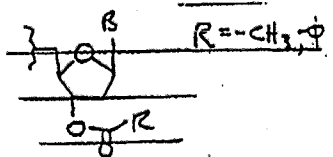can be used in the first step.

Photolysis of the photolabile group NV or NVOC on the
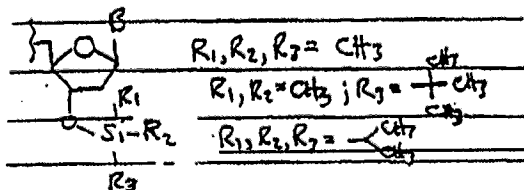5' positions of the nucleosides is carried out at ~362 nm with

72.

72

an intensity of 14 mW/cm$^2$ for 10 minutes with the substrate
side (side containing the photolabile group) immersed in
dioxane.  After the coupling of the next nucleoside is
complete, the photolysis is repeated followed by another
coupling until the desired oligomer is obtained.

One of the most common 3'-O-protecting groups is the
ester, in particular the acetate:

R = -CH$_3$, Φ

The groups can be removed by mild base treatment 0.1N
NaOH/MeOH or K$_2$CO$_3$/H$_2$O/MeOH.

Another group used most often is the silyl ether:

$R_1, R_2, R_3 = CH_3$

$R_1, R_2 = CH_3$ ; $R_3 = $

$R_1, R_2, R_3 = $

These groups can be removed by neutral conditions
using 1 M tetra-n-butylammonium fluoride in THF or under acid
conditions.

With respect to photodeprotection, the nitroveratryl
group could also be used to protect the 3'-position.

73

Here, light (photolysis) would be used to remove these protecting groups.

A variety of ethers can also be used in the protection of the 3'-O-position:



Removal of these groups usually involves acid or catalytic methods.

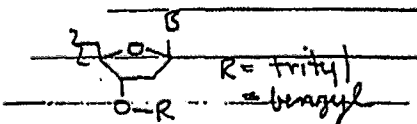Note that corresponding linkages and photoblocked amino acids are described in detail in Serial No. 07/624,120, now abandoned, which is hereby incorporated herein by reference.

Although the specificity of interactions at particular locations will usually be homogeneous due to a homogeneous polymer being synthesized at each defined location, for certain purposes, it may be useful to have mixed polymers with a commensurate mixed collection of interactions occurring at specific defined locations, or degeneracy reducing analogues, which have been discussed above and show broad specificity in binding. Then, a positive interaction signal may result from any of a number of sequences contained therein.

As an alternative method of generating a matrix pattern on a substrate, preformed polymers may be individually attached at particular sites on the substrate. This may be performed by individually attaching reagents one at a time to specific positions on the matrix, a process which may be automated. See, e.g., Serial No. 07/435,316, now abandoned, and Barrett et al. (1993) U.S. Pat. No. 5,252,743. Another way of generating a positionally defined matrix pattern on a substrate is to have individually specific reagents which interact with each specific position on the substrate. For example, oligonucleotides may be synthesized at defined locations on the substrate. Then the substrate would have on
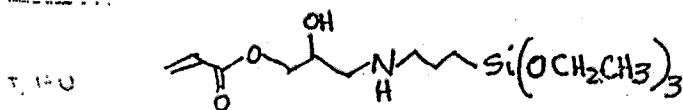
74

its surface a plurality of regions having homogeneous
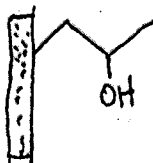oligonucleotides attached at each position.

In particular, at least four different substrate
preparation procedures are available for treating a substrate
surface.  They are the standard VLSIPS™ Technology method,
polymeric substrates, Durapore™, and synthetic beads or fibers.
The treatment labeled "standard VLSIPS™ Technology" method is
described in Serial No. 07/624,120, now abandoned, and involves
applying amino-propyltriethoxysilane to a glass surface.

The polymeric substrate approach involves either of
two ways of generating a polymeric substrate.  The first uses a
high concentration of aminopropyltriethoxysilane (2-20%) in an
aqueous ethanol solution (95%).  This allows the silane
compound to polymerize both in solution and on the substrate
surface, which provides a high density of amines on the surface
of the glass.  This density is contrasted with the standard
VLSIPS method.  This polymeric method allows for the deposition
on the substrate surface of a monolayer due to the anhydrous
method used with the aforementioned silane.

The second polymeric method involves either the
coating or covalent binding of an appropriate acrylic acid
polymer onto the substrate surface.  In particular, e.g., in
DNA synthesis, a monomer such as a hydroxypropylacrylate is
used to generate a high density of hydroxyl groups on the
substrate surface, allowing for the formation of phosphate
bonds.  An example of such a compound is shown:



The method using a Durapore™ membrane (Millipore)
consists of a polyvinylidine difluoride coating with
crosslinked polyhydroxylpropyl acrylate [PVDF-HPA]:

75

Here the building up of, e.g., a DNA oligomer, can be started
immediately since phosphate bonds to the surface can be
accomplished in the first step with no need for modification.
A nucleotide dimer (5'-C-T-3') has been successfully made on
this substrate.

The fourth method utilizes synthetic beads or fibers.
This would use another substrate, such as a teflon copolymer
graft bead or fiber, which is covalently coated with an organic
layer (hydrophilic) terminating in hydroxyl sites (commercially
available from Molecular Biosystems, Inc.) This would offer
the same advantage as the Durapore™ membrane, allowing for
immediate phosphate linkages, but would give additional contour
by the 3-dimensional growth of oligomers.

A matrix pattern of new reagents may be targeted to
each specific oligonucleotide position by attaching a
complementary oligonucleotide to which the substrate bound form
is complementary. For instance, a number of regions may have
homogeneous oligonucleotides synthesized at various locations.
Oligonucleotide sequences complementary to each of these can
be individually generated and linked to a particular specific
reagents. Often these specific reagents will be antibodies.
As each of these is specific for finding its complementary
oligonucleotide, each of the specific reagents will bind
through the oligonucleotide to the appropriate matrix position.
A single step having a combination of different specific
reagents being attached specifically to a particular
oligonucleotide will thereby bind to its complement at the
defined matrix position. The oligonucleotides will typically
then be covalently attached, using, e.g., an acridine dye, for
photocrosslinking. Psoralen is a commonly used acridine dye
for photocrosslinking purposes, see, e.g., Song et al. (1979)
Photochem. Photobiol. 29:1177-1197; Cimino et al. (1985) Ann.
Rev. Biochem. 54:1151-1193; Parsons (1980) Photochem.
Photobiol. 32:813-821; and Dattagupta et al. (1985) U.S. Pat.
No. 4,542,102, and (1987) U.S. Pat. No. 4,713,326; each of
which is hereby incorporated herein by reference. This method

76

allows a single attachment manipulation to attach all of the
specific reagents to the matrix at defined positions and
results in the specific reagents being homogeneously located at
defined positions. In many embodiments, the specific reagents
will be antibodies.

In an alternative embodiment, antibody molecules may
be used to specifically direct binding to defined positions on
a substrate. The VLSIPS technology may be used to generate
specific epitopes at each position on the substrate. Antibody
molecules having specificity of interaction may be used to
attach oligonucleotides, thereby avoiding the interference of
internal polynucleotide sequences from binding to the substrate
complementary oligonucleotides. In fact, the specificity of
interaction for positional targeting may be achieved by use of
nucleotide analogues which do not interact with the natural
nucleotides. For example, other synthetic nucleotides have
been made which undergo base pairing, thereby providing the
specificity of targeting, but the synthetic nucleotides also do
not interact with the natural biological nucleotides. Thus,
synthetic oligonucleotides would be useful for attachment to
biological nucleotides and specific targeting. Moreover, the
VLSIPS synthetic processes would be useful in generating the
VLSIPS substrate, and standard oligonucleotide synthesis could
be applied, with minor modifications, to produce the
complementary sequences which would be attached to other
specific reagents.

D.   Surface Immobilization
     1.   caged biotin
An alternative method of attaching reagents in a
positionally defined matrix pattern is to use a caged biotin
system. See Barrett et al. (1993) U.S. Pat. No. 5,252,743,
which is hereby incorporated herein by reference, for
additional details on the chemistry and application of caged
biotin embodiments. In short, the caged biotin has a
photosensitive blocking moiety which prevents the combination
of avidin to biotin. At positions where the photo-lithographic
process has removed the blocking group, high affinity biotin

᠁᠁

77

sites are generated.  Thus, by a sequential series of
photolithographic deblocking steps interspersed with exposure
of those regions to appropriate biotin containing reagents,
only those locations where the deblocking takes place will form
an avidin-biotin interaction.  Because the avidin-biotin
binding is very tight, this will usually be virtually
irreversible binding."

       2.    crosslinked interactions
         The surface immobilization may also take place by
photo crosslinking of defined oligonucleotides linked to
specific reagents.  After hybridization of the complementary
oligonucleotides, the oligonucleotides may be crosslinked by a
reagent by psoralen or another similar type of acridine dye.
Other useful cross linking reagents are described in Dattagupta
et al. (1985) U.S. Pat. No. 4,542,102, and (1987) U.S. Pat. No.
4,713,326.
         In another embodiment, colony or phage plaque
transfer of biological polymers may be transferred directly
onto a silicon substrate.  For example, a colony plate may be
transferred onto a substrate having a generic oligonucleotide
sequence which hybridizes to another generic complementary
sequence contained on all of the vectors into which inserts are
cloned.  This will specifically only bind those molecules which
are actually contained in the vectors containing the desired
complementary sequence.  This immobilization allows for
producing a matrix onto which a sequence specific reagent can
bind, or for other purposes.  In a further embodiment, a
plurality of different vectors each having a specific
oligonucleotide attached to the vector may be specifically
attached to particular regions on a matrix having a
complementary oligonucleotide attached thereto.

78

VIII.    HYBRIDIZATION/SPECIFIC INTERACTION
    A.    General

        As discussed previously in the VLSIPS™ Technology
parent applications, the VLSIPS™ technology substrates may be
used for screening for specific interactions with sequence
specific targets or probes.

        In addition, the availability of substrates having
the entire repertoire of possible sequences of a defined length
opens up the possibility of sequencing by hybridization.  This
sequence may be de novo determination of an unknown sequence,
particularly of nucleic acid, verification of a sequence
determined by another method, or an investigation of changes in
a previously sequenced gene, locating and identifying specific
changes.  For example, often Maxam and Gilbert sequencing
techniques are applied to sequences which have been determined
by Sanger and Coulson.  Each of those sequencing technologies
have problems with resolving particular types of sequences.
Sequencing by hybridization may serve as a third and
independent method for verifying other sequencing techniques.
See, e.g., (1988) Science 242:1245.

        In addition, the ability to provide a large
repertoire of particular sequences allows use of short
subsequences and hybridization as a means to fingerprint a
sample.  This may be used in a nucleic acid, as well as other
polymer embodiments.  For example, fingerprinting to a high
degree of specificity of sequence matching may be used for
identifying highly similar samples, e.g., those exhibiting high
homology to the selected probes.  This may provide a means for
determining classifications of particular sequences.  This·
should allow determination of whether particular genomes of
bacteria, phage, or even higher cells might be related to one
another.

        In addition, fingerprinting may be used to identify
an individual source of biological sample.  See, e.g., Lander,
E. (1989) Nature, 339:501-505, and references therein.  For
example, a DNA fingerprint may be used to determine whether a
genetic sample arose from another individual.  This would be
particularly useful in various sorts of forensic tests to

79

determine, e.g., paternity or sources of blood samples.
Significant detail on the particulars of genetic fingerprinting
for identification purposes are described in, e.g., Morris et
al. (1989) "Biostatistical evolution of evidence from
continuous allele frequency distribution DNA probes in
reference to disputed paternity of identity," J. Forensic
Science 34:1311-1317; and Neufeld et al. (1990) Scientific
American 262:46-53; each of which is hereby incorporated herein
by reference.

In another embodiment, a fingerprinting-like
procedure may be used for classifying cell types by analyzing a
pattern of specific nucleic acids present in the cell.  A
series of antibodies may be used to identify cell markers,
e.g., proteins, usually on the cell surface, but intracellular
markers may also be used.  Antigens which are extracellularly
expressed are preferred so cell lysis is unnecessary in the
screening, but intracellular markers may also be useful.  The
markers will usually be proteins, but may be nucleic acids,
lipids, metabolites, carbohydrates, or other cellular
components.  See, e.g., Winkelgren, I. (1990) Science News
136:234-237, which indicates extracellular DNA may be common,
and suggesting that such might be characteristic of cell types,
stage, or physiology.  This may also be useful in defining the
temporal stage of development of cells, e.g., stem cells or
other cells which undergo temporal changes in development.  For
example, the stage of a cell, or group of cells, may be tested
or defined by isolating a sample of mRNA from the population
and testing to see what sequences are present in messenger
populations.  Direct samples, or amplified samples, may be
used.  Where particular mRNA or other nucleic acid sequences
may be characteristic of or shown to be characteristic of
particular developmental stages, physiological states, or other
conditions, this fingerprinting method may define them.
Similar sorts of fingerprinting may be used for determining T-
cell classes or perhaps even to generate classification schemes
for such proteins as major histocompatibility complex antigens.
Thus, the ability to make these substrates allows both the
generation of reagents which will be used for defining

80

subclasses or classes of cells or other biological materials, but also provides the mechanisms for selecting those cells which may be found in defined population groups.

In addition to cell classification defined by such a combination of properties, typically expression of extracellular antigens, the present invention also provides the means for isolating homogeneous population of cells. Once the antigenic determinants which define a cell class have been identified, these antigens may be used in a sequential selection process to isolate only those cells which exhibit the combination of defining structural properties.

The present invention may also be used for mapping sequences within a larger segment. This may be performed by at least two methods, particularly in reference to nucleic acids. Often, enormous segments of DNA are subcloned into a large plurality of subsequences. Ordering these subsequences may be important in determining the overlaps of sequences upon nucleotide determinations. Mapping may be performed by immobilizing particularly large segments onto a matrix using the VLSIPS™ Technology. Alternatively, sequences may be ordered by virtue of subsequences shared by overlapping segments. See, e.g., Craig et al. (1990) Nuc. Acids Res. 18:2653-2660; Michiels et al. (1987) CABIOS 3:203-210; and Olson et al. (1986) Proc. Natl. Acad. Sci. USA 83:7826-7830.

B.    Important Parameters

The extent of specific interaction between reagents immobilized to the VLSIPS™ Technology substrate and another sequence specific reagent may be modified by the conditions of the interaction. Sequencing embodiments typically require high fidelity hybridization and the ability to discriminate perfect matching from imperfect matching. Fingerprinting and mapping embodiments may be performed using less stringent conditions, depending upon the circumstances.

For example, the specificity of antibody/antigen interaction may depend upon such parameters as pH, salt concentration, ionic composition, solvent composition, detergent composition and concentration, and chaotropic agent

81

concentration.  See, e.g., Harlow and Lane (1988) Antibodies: A Laboratory Manual, Cold Spring Harbor Press, New York.  By careful control of these parameters, the affinity of binding may be mapped across different sequences.

In a nucleic acid hybridization embodiment, the specificity and kinetics of hybridization have been described in detail by, e.g., Wetmur and Davidson (1968) J. Mol. Biol., 31:349-370, Britten and Kohne (1968) Science 161:529-530, and Kanehisa, (1984) Nuc. Acids Res. 12:203-213, each of which is hereby incorporated herein by reference.  Parameters which are well known to affect specificity and kinetics of reaction include salt conditions, ionic composition of the solvent, hybridization temperature, length of oligonucleotide matching sequences, guanine and cytosine (GC) content, presence of hybridization accelerators, pH, specific bases found in the matching sequences, solvent conditions, and addition of organic solvents.

In particular, the salt conditions required for driving highly mismatched sequences to completion typically include a high salt concentration.  The typical salt used is sodium chloride (NaCl), however, other ionic salts may be utilized, e.g., KCl.  Depending on the desired stringency hybridization, the salt concentration will often be less than about 3 molar, more often less than 2.5 molar, usually less than about 2 molar, and more usually less than about 1.5 molar. For applications directed towards higher stringency matching, the salt concentrations would typically be lower.  Ordinary high stringency conditions will utilize salt concentration of less than about 1 molar, more often less then about 750 millimolar, usually less than about 500 millimolar, and may be as low as about 250 or 150 millimolar.

The kinetics of hybridization and the stringency of hybridization both depend upon the temperature at which the hybridization is performed and the temperature at which the washing steps are performed.  Temperatures at which steps for low stringency hybridization are desired would typically be lower temperatures, e.g., ordinarily at least about 15°C, more ordinarily at least about 20°C, usually at least about 25°C, and

82

more usually at least about 30°C. For those applications
requiring high stringency hybridization, or fidelity of
hybridization and sequence matching, temperatures at which
hybridization and washing steps are performed would typically
be high. For example, temperatures in excess of about 35°C
would often be used, more often in excess of about 40°C,
usually at least about 45°C, and occasionally even temperatures
as high as about 50°C or 60°C or more. Of course, the
hybridization of oligonucleotides may be disrupted by even
higher temperatures. Thus, for stripping of targets from
substrates, as discussed below, temperatures as high as 80°C,
or even higher may be used.

The base composition of the specific oligonucleotides
involved in hybridization affects the temperature of melting,
and the stability of hybridization as discussed in the above
references. However, the bias of GC rich sequences to
hybridize faster and retain stability at higher temperatures
can be compensated for by the inclusion in the hybridization
incubation or wash steps of various buffers. Sample buffers
which accomplish this result include the triethly-and trimethyl
ammonium buffers. See, e.g., Wood et al. (1987) Proc. Natl.
Acad. Sci. USA, 82:1585-1588, and Khrapko, K. et al. (1989)
FEBS Letters 256:118-122.

The rate of hybridization can also be affected by the
inclusion of particular hybridization accelerators. These
hybridization accelerators include the volume exclusion agents
characterized by dextran sulfate, or polyethylene glycol (PEG).
Dextran sulfate is typically included at a concentration of
between 1% and 40% by weight. The actual concentration
selected depends upon the application, but typically a faster
hybridization is desired in which the concentration is
optimized for the system in question. Dextran sulfate is often
included at a concentration of between 0.5% and 2% by weight or
dextran sulfate at a concentration between about 0.5% and 5%.
Alternatively, proteins which accelerate hybridization may be
added, e.g., the recA protein found in E. coli or other
homologous proteins.

83

With respect to those embodiments where specific
reagents are not oligonucleotides, the conditions of specific
interaction would depend on the affinity of binding between the
specific reagent and its target. Typically parameters which
would be of particular importance would be pH, salt
concentration anion and cation compositions, buffer
concentration, organic solvent inclusion, detergent
concentration, and inclusion of such reagents such as
chaotropic agents. In particular, the affinity of binding may
be tested over a variety of conditions by multiple washes and
repeat scans or by using reagents with differences in binding
affinity to determine which reagents bind or do not bind under
the selected binding and washing conditions. The spectrum of
binding affinities may provide an additional dimension of
information which may be very useful in identification purposes
and mapping.

Of course, the specific hybridization conditions will
be selected to correspond to a discriminatory condition which
provides a positive signal where desired but fails to show a
positive signal at affinities where interaction is not desired.
This may be determined by a number of titration steps or with
a number of controls which will be run during the hybridization
and/or washing steps to determine at what point the
hybridization conditions have reached the stage of desired
specificity.

IX.   DETECTION METHODS

Methods for detection depend upon the label selected.
The criteria for selecting an appropriate label are discussed
below, however, a fluorescent label is preferred because of its
extreme sensitivity and simplicity. Standard labeling
procedures are used to determine the positions where
interactions between a sequence and a reagent take place. For
example, if a target sequence is labeled and exposed to a
matrix of different probes, only those locations where probes
do interact with the target will exhibit any signal.
Alternatively, other methods may be used to scan the matrix to
determine where interaction takes place. Of course, the

84

spectrum of interactions may be determined in a temporal manner
by repeated scans of interactions which occur at each of a
multiplicity of conditions. However, instead of testing each
individual interaction separately, a multiplicity of sequence
interactions may be simultaneously determined on a matrix.

A.    Labeling Techniques

The target polynucleotide may be labeled by any of a
number of convenient detectable markers. A fluorescent label
is preferred because it provides a very strong signal with low
background. It is also optically detectable at high resolution
and sensitivity through a quick scanning procedure. Other
potential labeling moieties include, radioisotopes,
chemiluminescent compounds, labeled binding proteins, heavy
metal atoms, spectroscopic markers, magnetic labels, and linked
enzymes.

Another method for labeling may bypass any label of
the target sequence. The target may be exposed to the probes,
and a double strand hybrid is formed at those positions only.
Addition of a double strand specific reagent will detect where
hybridization takes place. An intercalative dye such as
ethidium bromide may be used as long as the probes themselves
do not fold back on themselves to a significant extent forming
.hairpin loops. See, e.g., Sheldon et al. (1986) U.S. Pat. No.
4,582,789. However, the length of the hairpin loops in short
oligonucleotide probes would typically be insufficient to form
a stable duplex.

In another embodiment, different targets may be
simultaneously sequenced where each target has a different
label. For instance, one target could have a green fluorescent
label and a second target could have a red fluorescent label.
The scanning step will distinguish sites of binding of the red
label from those binding the green fluorescent label. Each
sequence can be analyzed independently from one another.

Suitable chromogens will include molecules and
compounds which absorb light in a distinctive range of
wavelengths so that a color may be observed, or emit light when
irradiated with radiation of a particular wave length or wave

85

length range, e.g., fluorescers. Biliproteins, e.g.,
phycoerythrin, may also serve as labels.

A wide variety of suitable dyes are available, being
primarily chosen to provide an intense color with minimal
absorption by their surroundings. Illustrative dye types
include quinoline dyes, triarylmethane dyes, acridine dyes,
alizarine dyes, phthaleins, insect dyes, azo dyes,
anthraquinoid dyes, cyanine dyes, phenazathionium dyes, and
phenazoxonium dyes.

A wide variety of fluorescers may be employed either
by themselves or in conjunction with quencher molecules.
Fluorescers of interest fall into a variety of categories
having certain primary functionalities. These primary
functionalities include 1- and 2-aminonaphthalene, p,p'-
diaminostilbenes, pyrenes, quaternary phenanthridine salts, 9-
aminoacridines, p,p'-diaminobenzophenone imines, anthracenes,
oxacarbocyanine, merocyanine, 3-aminoequilenin, perylene, bis-
benzoxazole, bis-p-oxazolyl benzene, 1,2-benzophenazin,
retinol, bis-3-aminopyridinium salts, hellebrigenin,
tetracycline, sterophenol, benzimidzaolylphenylamine, 2-oxo-3-
chromen, indole, xanthen, 7-hydroxycoumarin, phenoxazine,
salicylate, strophanthidin, porphyrins, triarylmethanes and
flavin. Individual fluorescent compounds which have
functionalities for linking or which can be modified to
incorporate such functionalities include, e.g., dansyl
chloride; fluoresceins such as 3,6-dihydroxy-9-
phenylxanthhydrol; rhodamineisothiocyanate; N-phenyl 1-amino-8-
sulfonatonaphthalene; N-phenyl 2-amino-6-sulfonatonaphthalene;
4-acetamido-4-isothiocyanato-stilbene-2,2'-disulfonic acid;
pyrene-3-sulfonic acid; 2-toluidinonaphthalene-6-sulfonate; N-
phenyl, N-methyl 2-aminoaphthalene-6-sulfonate; ethidium
bromide; stebrine; auromine-0,2-(9'-anthroyl)palmitate; dansyl
phosphatidylethanolamine; N,N'-dioctadecyl oxacarbocyanine;
N,N'-dihexyl oxacarbocyanine; merocyanine, 4-
(3'pyrenyl)butyrate; d-3-aminodesoxy-equilenin; 12-(9'-
anthroyl)stearate; 2-methylanthracene; 9-vinylanthracene; 2,2'-
(vinylene-p-phenylene)bisbenzoxazole; p-bis[2-(4-methyl-5-
phenyl-oxazolyl)]benzene; 6-dimethylamino-1,2-benzophenazin;

86

retinol; bis(3'-aminopyridinium) 1,10-decandiyl diiodide; sulfonaphthylhydrazone of hellibrienin; chlorotetracycline; N-(7-dimethylamino-4-methyl-2-oxo-3-chromenyl)maleimide; N-[p-(2-benzimidazolyl)-phenyl]maleimide; N-(4-fluoranthyl)maleimide; bis(homovanillic acid); resazarin; 4-chloro-7-nitro-2,1,3-benzooxadiazole; merocyanine 540; resorufin; rose bengal; and 2,4-diphenyl-3(2H)-furanone.

Desirably, fluorescers should absorb light above about 300 nm, preferably about 350 nm, and more preferably above about 400 nm, usually emitting at wavelengths greater than about 10 nm higher than the wavelength of the light absorbed. It should be noted that the absorption and emission characteristics of the bound dye may differ from the unbound dye. Therefore, when referring to the various wavelength ranges and characteristics of the dyes, it is intended to indicate the dyes as employed and not the dye which is unconjugated and characterized in an arbitrary solvent.

Fluorescers are generally preferred because by irradiating a fluorescer with light, one can obtain a plurality of emissions. Thus, a single label can provide for a plurality of measurable events.

Detectable signal may also be provided by chemiluminescent and bioluminescent sources. Chemiluminescent sources include a compound which becomes electronically excited by a chemical reaction and may then emit light which serves as the detectible signal or donates energy to a fluorescent acceptor. A diverse number of families of compounds have been found to provide chemiluminescence under a variety of conditions. One family of compounds is 2,3-dihydro-1,-4-phthalazinedione. The most popular compound is luminol, which is the 5-amino compound. Other members of the family include the 5-amino-6,7,8-trimethoxy- and the dimethylamino[ca]benz analog. These compounds can be made to luminesce with alkaline hydrogen peroxide or calcium hypochlorite and base. Another family of compounds is the 2,4,5-triphenylimidazoles, with lophine as the common name for the parent product. Chemiluminescent analogs include para-dimethylamino and -methoxy substituents. Chemiluminescence may also be obtained

87

87

with oxalates, usually oxalyl active esters, e.g., p-nitrophenyl and a peroxide, e.g., hydrogen peroxide, under basic conditions. Alternatively, luciferins may be used in conjunction with luciferase or lucigenins to provide bioluminescence.

Spin labels are provided by reporter molecules with an unpaired electron spin which can be detected by electron spin resonance (ESR) spectroscopy. Exemplary spin labels include organic free radicals, transitional metal complexes, particularly vanadium, copper, iron, and manganese, and the like. Exemplary spin labels include nitroxide free radicals.

B.    Scanning System

With the automated detection apparatus, the correlation of specific positional labeling is converted to the presence on the target of sequences for which the reagents have specificity of interaction. Thus, the positional information is directly converted to a database indicating what sequence interactions have occurred. For example, in a nucleic acid hybridization application, the sequences which have interacted between the substrate matrix and the target molecule can be directly listed from the positional information. The detection system used is described in Pirrung et al. (1992) U.S. Pat. No. 5,143,854; and Serial No. 07/624,120, now abandoned. Although the detection described therein is a fluorescence detector, the detector may be replaced by a spectroscopic or other detector. The scanning system may make use of a moving detector relative to a fixed substrate, a fixed detector with a moving substrate, or a combination. Alternatively, mirrors or other apparatus can be used to transfer the signal directly to the detector. See, e.g, Serial No. 07/624,120, now abandoned, which is hereby incorporated herein by reference.

The detection method will typically also incorporate some signal processing to determine whether the signal at a particular matrix position is a true positive or may be a spurious signal. For example, a signal from a region which has actual positive signal may tend to spread over and provide a positive signal in an adjacent region which actually should not

88

have one. This may occur, e.g., where the scanning system is not properly discriminating with sufficiently high resolution in its pixel density to separate the two regions. Thus, the signal over the spatial region may be evaluated pixel by pixel to determine the locations and the actual extent of positive signal. A true positive signal should, in theory, show a uniform signal at each pixel location. Thus, processing by plotting number of pixels with actual signal intensity should have a clearly uniform signal intensity. Regions where the signal intensities show a fairly wide dispersion, may be particularly suspect and the scanning system may be programmed to more carefully scan those positions.

In another embodiment, as the sequence of a target is determined at a particular location, the overlap for the sequence would necessarily have a known sequence. Thus, the system can compare the possibilities for the next adjacent position and look at these in comparison with each other. Typically, only one of the possible adjacent sequences should give a positive signal and the system might be programmed to compare each of these possibilities and select that one which gives a strong positive. In this way, the system can also simultaneously provide some means of measuring the reliability of the determination by indicating what the average signal to background ratio actually is.

More sophisticated signal processing techniques can be applied to the initial determination of whether a positive signal exists or not. See, e.g., Serial No. 07/624,120, now abandoned.

From a listing of those sequences which interact, data analysis may be performed on a series of sequences. For example, in a nucleic acid sequence application, each of the sequences may be analyzed for their overlap regions and the original target sequence may be reconstructed from the collection of specific subsequences obtained therein. Other sorts of analyses for different applications may also be performed, and because the scanning system directly interfaces with a computer the information need not be transferred manually. This provides for the ability to handle large

89

amounts of data with very little human intervention. This, of
course, provides significant advantages over manual
manipulations. Increased throughput and reproducibility is
thereby provided by the automation of a vast majority of steps
in any of these applications.

XI.    DATA ANALYSIS
    A.    General
        Data analysis will typically involve aligning the
proper sequences with their overlaps to determine the target
sequence. Although the target "sequence" may not specifically
correspond to any specific molecule, especially where the
target sequence is broken and fragmented in the sequencing
process, the sequence corresponds to a contiguous sequence of
the subfragments.

        The data analysis can be performed by a computer
using an appropriate program. See, e.g., Drmanac, R. et al.
(1989) Genomics 4:114-128; and a commercially available
analysis program available from the Genetic Engineering Center,
P.O. Box 794, 11000 Belgrade, Yugoslavia. Although the
specific manipulations necessary to reassemble the target
sequence from fragments may take many forms, one embodiment
uses a sorting program to sort all of the subsequences using a
defined hierarchy. The hierarchy need not necessarily
correspond to any physical hierarchy, but provides a means to
determine, in order, which subfragments have actually been
found in the target sequence. In this manner, overlaps can be
checked and found directly rather than having to search
throughout the entire set after each selection process. For
example, where the oligonucleotide probes are 10-mers, the
first 9 positions can be sorted. A particular subsequence can
be selected as in the examples, to determine where the process
starts. As analogous to the theoretical example provided
above, the sorting procedure provides the ability to
immediately find the position of the subsequence which contains
the first 9 positions and can compare whether there exists more
than 1 subsequence during the first 9 positions. In fact, the
computer can easily generate all of the possible target

'\()

90

sequences which contain given combination of subsequences.
Typically there will be only one, but in various situations,
there will be more.

An exemplary flow chart for a sequencing program is
provided in Figure ✗.  In general terms, the program provides
for automated scanning of the substrate to determine the
positions of probe and target interaction.  Simple processing
of the intensity of the signal may be incorporated to filter
out clearly spurious signals.  The positions with positive
interaction are correlated with the sequence specificity of
specific matrix positions, to generate the set of matching
subsequences.  This information is further correlated with
other target sequence information, e.g., restriction fragment
analysis.  The sequences are then aligned using overlap data,
thereby leading to possible corresponding target sequences
which will, optimally, correspond to a single target sequence.

      B.    Hardware

      A variety of computer systems may be used to run a
sequencing program.  The program may be written to provide both
the detecting and scanning steps together and will typically be
dedicated to a particular scanning apparatus.  However, the
components and functional steps may be separated and the
scanning system may provide an output, e.g., through tape or an
electronic connection into a separate computer which separately
runs the sequencing analysis program.  The computer may be any
of a number of machines provided by standard computer
manufacturers, e.g., IBM compatible machines, Apple™ machines,
VAX machines, and others, which may often use a UNIX™ operating
system.  Of course, the hardware used to run the analysis
program will typically determine what programming language
would be used.

      C.    Software

      Software would be easily developed by a person of
ordinary skill in the programming art, following the flow chart
provided, or based upon the input provided and the desired
result.

91

Of course, an exemplary embodiment is a
polynucleotide sequence system. However, the theoretical and
mathematical manipulations necessary for data analysis of other
linear molecules, such as polypeptides, carbohydrates, and
various other polymers are conceptually similar. Simple
branching polymers will usually also be sequencable using
similar technology. However, where there is branching, it may
be desired that additional recognition reagents be used to
determine the nature and location of branches. This can easily
be provided by use of appropriate specific reagents which would
be generated by methods similar to those used to produce
specific reagents for linear polymers.

XII. SUBSTRATE REUSE

Where a substrate is made with specific reagents that
are relatively insensitive to the handling and processing steps
involved in a single cycle of use, the substrate may often be
reused. The target molecules are usually stripped off of the
solid phase specific recognition molecules. Of course, it is
preferred that the manipulations and conditions be selected as
to be mild and to not affect the substrate. For example, if a
substrate is acid labile, a neutral pH would be preferred in
all handling steps. Similar sensitivities would be carefully
respected where recycling is desired.

A.    Removal of Label

Typically for a recycling, the previously attached
specific interaction would be disrupted and removed. This will
typically involve exposing the substrate to conditions under
which the interaction between probe and target is disrupted.
Alternatively, it may be exposed to conditions where the target
is destroyed. For example, where the probes are
oligonucleotides and the target is a polynucleotide, a heating
and low salt wash will often be sufficient to disrupt the
interactions. Additional reagents may be added such as
detergents, and organic or inorganic solvents which disrupt the
interaction between the specific reagents and target. In an
embodiment where the specific reagents are antibodies, the

92

92

substrate may be exposed to a gentle detergent which will denature the specific binding between the antibody and its target. The conditions are selected to avoid severe disruption or destruction of the structure of the antibody and to maintain the specificity of the antibody binding site. Conditions with specific pH, detergent concentration, salt concentration, ionic concentration, and other parameters may be selected which disrupt the specific interactions.

B.    Storage and Preservation

As indicated above, the matrix will typically be maintained under conditions where the matrix itself and the linkages and specific reagents are preserved. Various specific preservatives may be added which prevent degradation. For example, if the reagents are acid or base labile, a neutral pH buffer will typically be added. It is also desired to avoid destruction of the matrix by growth of organisms which may destroy organic reagents attached thereto. For this reason, a preservative such as cyanide or azide may be added. However, the chemical preservative should also be selected to preserve the chemical nature of the linkages and other components of the substrate. Typically, a detergent may also be included.

C.    Processes to Avoid Degradation of Oligomers

In particular, a substrate comprising a large number of oligomers will be treated in a fashion which is known to maintain the quality and integrity of oligonucleotides. These include storing the substrate in a carefully controlled environment under conditions of lower temperature, cation depletion (EDTA and EGTA), sterile conditions, and inert argon or nitrogen atmosphere.

XIII.    INTEGRATED SEQUENCING STRATEGY

A.    Initial Mapping Strategy

As indicated above, although the VLSIPS™ technology may be applied to sequencing embodiments, it is often useful to integrate other concepts to simplify the sequencing. For example, nucleic acids may be easily sequenced by careful

93

selection of the vectors and hosts used for amplifying and
generating the specific target sequences. For example, it may
be desired to use specific vectors which have been designed to
interact most efficiently with the VLSIPS substrate. This is
also important in fingerprinting and mapping strategies. For
example, vectors may be carefully selected having particular
complementary sequences which are designed to attach to a
genetic or specific oligomer on the substrate. This is also
applicable to situations where it is desired to target
particular sequences to specific locations on the matrix.

In one embodiment, unnatural oligomers may be used to
target natural probes to specific locations on the VLSIPS
substrate. In addition, particular probes may be generated for
the mapping embodiment which are designed to have specific
combinations of characteristics. For example, the construction
of a mapping substrate may depend upon use of another automated
apparatus which takes clones isolated from a chromosome walk
and attaches them individually or in bulk to the VLSIPS
substrate.

In another embodiment, a variety of specific vectors
having known and particular "targeting" sequences adjacent to
the cloning sites may be individually used to clone a selected
probe, and the isolated probe will then be targetable to a site
on the VLSIPS substrate with a sequence complementary to the
"target" sequence.

    B.    <u>Selection of Smaller Clones</u>

In the fingerprinting and mapping embodiments, the
selection of probes may be very important. Significant
mathematical analysis may be applied to determine which
specific sequences should be used as those probes. Of course,
for fingerprinting use, these sequences would be most desired
that show significant heterogeneity across the human
population. Selection of the specific sequences which would
most favorably be utilized will tend to be single copy
sequences within the genome.

Various hybridization selection procedures may be
applied to select sequences which tend not to be repeated

14